# LA-UR-13-26006

Title:        Functional and Performance Assessment of Erasure
              Coded Storage Systems

Author(s):    Sanchez, Taylor E
              Sackos, Joshua P
              Crossman, Blair A

Intended for:  technical report

Issued:        2013-07-30

## Los Alamos
### NATIONAL LABORATORY
—— EST. 1943 ——

# Functional and Performance Assessment of Erasure Coded Storage Systems

Taylor Sanchez - California State University San Bernardino

Josh Sackos - Washington State University

Blair Crossman - New Mexico Tech

Mentors: HB Chen and Jeff Inman

July 30, 2013

# Contents

# 1 Abstract

Current storage systems using RAID technologies were not designed to support exascale computing systems[4]. Multiple disk drives failures (more than two) have not been an issue in the past. Conventional hard disk drive capacities continue to grow beyond the terabyte range and storage systems continue to grow to hundreds of petabytes and even exabytes. The need to efficiently handle multiple disk drive failures is now a reality and a challenging problem. Rather than separating the data from error correction or check data as parity and CRCs do, erasure codes expand the data, adding redundancy so even if a portion of the data is mangled or lost, the original data can be retrieved from the remaining portion.[4] In recent years, various research and development efforts have attempted to create an erasure code and disk-based archive storage system with power-efficiency, data integrity, and system reliability[8]. In this project, we plan to conduct function and performance assessments of erasure code storage systems. We plan to build two small test beds for an all disk parallel archive storage system and conduct tests erasure code storage software from Scality[7] and Caringo[2].

As we approach exascale computing, disk storage becomes an attractive option due to scalability of disk bandwidth over tape drive[3]. Disk failure rates make exascale archive systems prone to data loss. Replication is a possible solution, but can double or triple required storage space. Erasure coding is a promising option for an exascale archive because it offers the durability of replication with less overhead[5]. The functionality of an erasure code archive system remains untested at exascale. Our project was to build and verify the functionality of two prototype erasure storage archives using commercial products from Scality and Caringo.

Both products had the functionality to read, write, balance, and rebuild data as well as offering metadata access. Caringo did not provide us with a POSIX gateway, but has a metadata indexing tool that allowed querying. We did not have the Scality indexing tool to query the metadata, but we were provided with the POSIX interface SFUSE. The POSIX gateway caused an average bandwidth loss of 70% for small files (less than 1MB) and 50% loss for large files (greater than 1GB).

# 2 Functional Tests

The following system setup and functional tests were performed.

1. Setup two commercial prototype systems for:

   (a) Scality
   (b) Caringo

2. Configured erasure coding software and built two archive storage systems.
3. Used Block storage device to access Object storage
4. Conducted the following tests on both systems:

   (a) Data ingesting (write operation)
   (b) Data retrieving (read operation)
   (c) Metadata querying (index and attribute searching)
   (d) Metadata accessing
   (e) Data load balance testing
   (f) Data repairing

5. Gathered preliminary data on POSIX interface overhead

# 3 Testbed

The servers used in this experiment consisted of 3 SuperMicro, 6 HP Proliant , and 10 IBM servers. The nodes were interconnected a via 1GigE switch. The switch was a 3COM SuperStack3 Switch 3870 24 port model. The terminal server used was a Cyclades-TS3000 terminal server.

SuperMicro SC745

- X7DWA-N Motherboard
- Dual Intel Xeon E5410 @ 2.33Ghz LGA771Processors
- 16GB and one server has 32GB (Modules: 4GB DDR2 @ 667Mhz)
- 1 WesternDigital 1TB HDD
- 1 Hitachi 4TB HDD (only on admin nodes)
- Redundant Power Supply

HP Proliant DL160 G6

- Dual Intel Xeon E5504 @ 2.00Ghz LGA1366 Processors
- 24GB (Modules: 2GB DDR3 @ 800Mhz)
- 32GB (Modules: 8GB DDR3 @800Mhz)
- 4 WesternDigital 1TB HDD

IBM System x3755

- Quad AMD 64 Opteron
- 32GB (Modules: 2GB DDR2)
- 4 WesternDigital 1TB HDD

## 3.1  Scality Testbed

A SuperMicro server with a 4TB Hitachi hard drive was used as the administration node. The storage nodes in the ring consisted of 6 HP Proliant servers, each with 3 x 1TB hard drives dedicated to ring storage space. Three of the storage nodes had 24GB of ram, and the other three 32 GB. The erasure coding scheme used in the Scality ring was n=3, and k=3.

## 3.2  Caringo Testbed

A SuperMicro server with 32 GB of RAM and a 4TB Hitachi hard drive was used as a metadata indexer. The storage nodes in the cluster consisted of 10 IBM servers, each with 4 x 1TB hard drives dedicated to storage space. When powered on, the storage nodes PXE boot a Caringo CAStor image into RAM. The erasure coding scheme used in the Caringo cluster was n=3, and k=3.

## 3.3  Network Overview

Figure  1 shows a graphical overview of the testbed network. The DHCP and DNS zone configuration files are listing  3 and listing  4 in appendix A. Both Caringo and Scality make use of round robin dns forwarding to load balance data streams.

Figure 1: Network Overview.



# 4 Scality

## 4.1 Overview

For every hard drive in the Scality ring there is a daemon running that is responsible for managing its respective disk. On each Proliant node there were 6 virtual machines that ran, making a total of 36 virtual machines in the Scality ring. Virtualization allows for a slightly more flexible setup for rebuilding data. Figure 4 in appendix B shows the Scality ring [6] that was implemented in this project. Every green oval is a virtual node.

## 4.2 Install

To install Scality in an existing Linux system, the proper services need to be installed on each node. Once installed Scality will need to be configured, see listing 5 and 6 in appendix C for details. Further customizations will require the sproxyd.conf file to be edited, see listing 7 in appendix C for details.

## 4.3 The Basics: Reading and Writing

The two interfaces for the ring in this project were the SFUSE (must be installed) and REST. Currently the RESTful interface limits the maximum size of a file to upload to 500MB. This limit will be removed future releases. See listing 11 in appendix F for details on SFUSE usage and installation.

### 4.3.1 RESTful
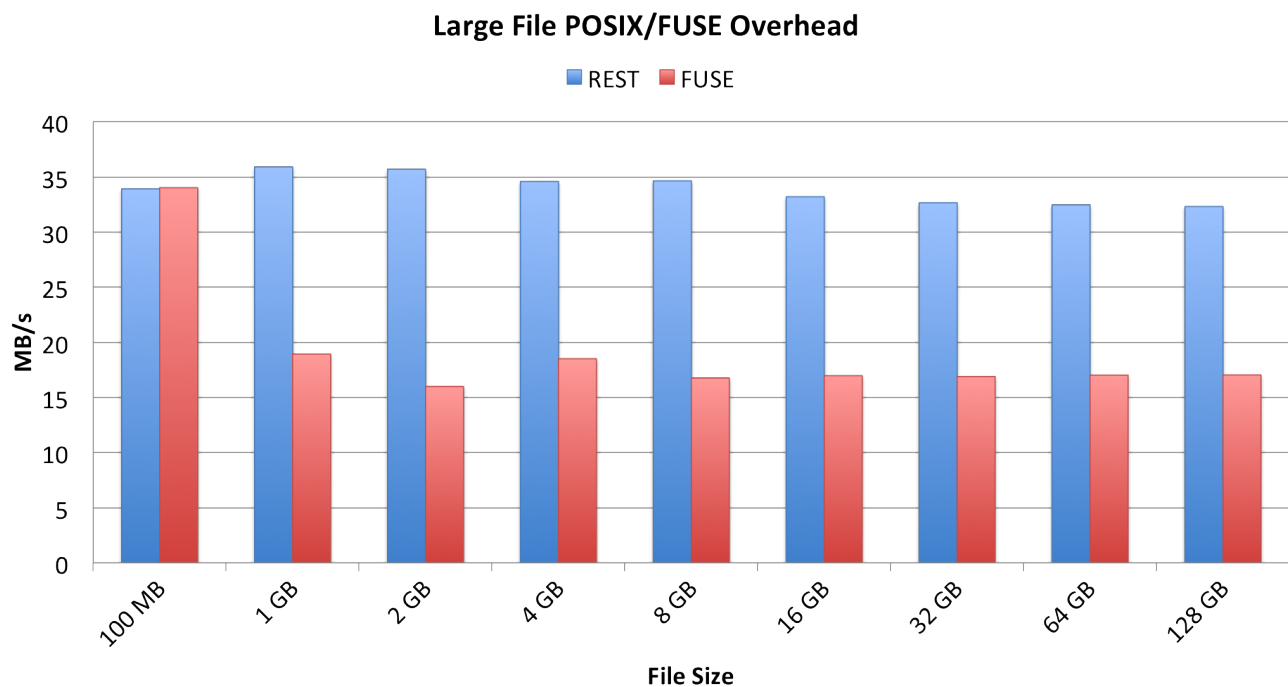
The service sproxyd will redirect your stream into the ring. Something to be aware of is that Curl loads the full file size into ram as it uploads the stream. This makes for a troublesome situation on low ram machines if Scality does remove their 500MB file upload limit. Appendix D contains a list of the curl commands.

### 4.3.2 POSIX

Scality provides a FUSE interface to communicate with the ring. While the POSIX interface is a requirement for the transition from tape storage to hard disk, it is less than ideal. Basic benchmarks of transferring different sized files through REST and through SFUSE (POSIX interface) interfaces were collected. In appendix F is the commands to run in bash to get SFUSE installed[6]. Once installed you should be able to use the simple cp command to put data on /ring/fuse/folder/file.name (notice it needs some sort of folder within the fuse directory, no files can go in the root of fuse). There is a large overhead with using the Scalitys POSIX interface, see figure 2 in for details.

There was not enough granularity in the tests to check where the breaking point is for the SFUSE's overhead. As you can see it does quite well on the 100MB files, but at 1GB it quickly drops. It seems reasonable to hypothesize that the 500MB data chunking size reduces the performance as SFUSE needs to start breaking up the data into chunks.

Figure 2: POSIX Overhead



### 4.4 Recovery Capabilities

Recovery requires a manual command to the supervisor. Data was recoverable during and after a rebuild, and when taking down servers. Performance overhead was not tested during a rebuild as fairly basic systems were used. Something worth mentioning about the supervisor is that it is not required for the ring to continue running, it can survive up to 24 hours without the supervisor available. A heartbeat fallback would probably be a good idea in a production system. Figure 3 is a screen shot where you can see the option to start the Repair action. Scality uses a type of virtual node that allows Scality to rebuild with less than required physical nodes. This means once you get the nodes back up, the ring

simply needs to rebalance the data to the new nodes.

Figure 3: Scality Supervisor



## 4.5 Data Balancing

This testing is very time consuming and was conducted throughout our testing. Although this was not the best documented we did see that the disks remained well balanced throughout our testing. The scality supervisor does list the storage usage on each node, as can be seen in the second screenshot of the Scality Overview section. One way to get a good simple look at the data balancing is to run 'df' on each node and check that they are all using a similar amount of each hard drive. The shortened output (listing 1) below was run after all testing was all finished and shows the disks to still be balanced through a rebuild, nodes being turned off for long time periods and having data written during the process.

Listing 1: Scality Balance

```
1  #Scality Balance
2  clush −w n[00,01,03,07,10,15] df 2> /dev/null | grep scality
3  #output:
4  #n00: /dev/sdb1 961432904 526904824 434528080 55% /scality/disk1
5  #n00: /dev/sdc1 961432904 524755448 436677456 55% /scality/disk2
6  #n00: /dev/sdd1 961432904 523089924 438342980 55% /scality/disk3
7  #n03: /dev/sdb1 961432904 527370188 434062716 55% /scality/disk1
8  #n03: /dev/sdc1 961432904 524312968 437119936 55% /scality/disk2
9  #n03: /dev/sdd1 961432904 522650476 438782428 55% /scality/disk3
10 #n07: /dev/sdb1 961432904 527218648 434214256 55% /scality/disk1
11 #...
```

7

## 4.6 Metadata

While Scality does not supply an indexing service out of the box, there are still ways of accessing the metadata contained in each server and file. The software Scality provides separately is called "Mesa". Without the Mesa software there are two ways of obtaining the metatdata in the system, one gets the metadata regarding the arc storage keys in the server, the other gets the metadata regarding the specific object.

A script to get the key-values from the servers is listed in appendix E, this produces a surprising amount of data.

The files will supply their metadata by using Curl with the corresponding simple command:

Listing 2: Single File MetaData

```
1    curl −0 −I http://localhost:81/proxy/chord/fileKey
2    curl −0 −I http://localhost:81/proxy/arc/fileKey
```

## 4.7 Scality Results and Observations

Scality didnt pass the metadata query requirement because it lacked the Mesa indexing service. Scality installs on top of an existing Linux operating system so it is very flexible and customizable. This stateful install is also inefficient as it takes up disk space that would otherwise be available to the storage cluster. The 500MB file size limit is small for the huge files that Los Alamos National Laboratory will need to push to their archive system. These problems seem possible to fix as Scality seems willing to work with their customers to improve their product. We ran a simple parallel write test on the ring to see if throughput would scale, and early testing indicates that it does.

# 5 Caringo

## 5.1 Overview

Caringo's CASTor software offers a bare-metal deployment. Each node is stateless, and can be provisioned either by usb key, or PXE boot from a network. CASTor nodes are decentralized. Each node replicates and encodes content, without any specialized nodes. When a new stream enters the CAStor cluster a "bid" is called from from the nodes in the cluster. Streams are then pushed to the node with the lowest bid, which is then responsible for encoding the stream and pushing out the data and code segments to other nodes in the cluster[2].

## 5.2 Install

The Caringo system installs via a RAM disk from a pxeboot setup. There is a config file that is kept on the admin node which will be delivered to the nodes during boot. When a new node comes online it needs access these files:

1. The castor file system image
2. The castor license file
3. The castor node configuration list: cfg-list.txt
4. The node configuration file pointed to by the list: node.cfg

See the examples in appendix G for how we set up DHCP for netbooting.

## 5.3 The Basics: Reading and Writing

Caringo's CAStor can be interfaced through either REST, or a commercially available Content File Server that allows data to be written to the cluster using a POSIX compliant NFS mount. The Content File server was not made available to us until after our testing deadline, so we were unable to access its functionality[2]. We were on the other had able to use the RESTful interface to interact with the CAStor cluster.

### 5.3.1 Basic Commands

Our testing was done through the RESTful interface utilizing the curl tool to create the http requests for us. Uploading objects to a CAStor cluster involves creating a bucket, and then storing objects in that bucket. The cluster will then return to you a GUID for that file, but you will still be able to access the file by simply performing a GET request on the URL that you stored the data in. One drawback of the RESTful interface is that streams are constantly moving between the nodes, so all RESTful requests must allow redirect when a 301 error is found.

See the command examples in appendix H.

### 5.3.2 POSIX

We received the Content File Server that serves as CAStor's POSIX interface after our testing phase had concluded[2]. We were unable to test its functionality. This interface was a requirement for the production setting as the current system is using a Block storage interface and will continue to require a Block storage interface until a transition to a RESTful backend is completed.

## 5.4 Recovery and Rebalance

A CAStor cluster is constantly rebuilding and rebalancing. In fact any volume or node left offline for over 14 days will be automatically reformatted and rebalanced. A supervisor may suspend rebuilding during maintenance, but the default behavior is to instantly begin a rebuild when a node goes offline. After removal of a node we found that all data could be successfully recovered, but if the number of nodes in the cluster fell beneath the total number of data and code segments, writing to the cluster would fail. For example in our cluster we had $3 + 3$ erasure code. If we had 6 nodes, we could write and read from the cluster. If we fell down to 5 nodes, we could read from the cluster, but no longer write to it until the cluster again had 6 nodes. The bidding system of the nodes allows for a near seamless rebuild while writing to the cluster, as new jobs will be assigned the nodes not responsible for the rebuild. As every node in the cluster is at once a supervisor node, compute node, and storage node, there is no single point of access to the cluster that could be tied up by a rebuild. This means that as long as you have enough nodes to compute a rebuild, regular cluster management and use is seamless. [2]

## 5.5 Metadata

Caringo's CAStor utilizes elastic search in order to index and query metadata. The indexer will index certain basic metadata fields for any file and offers the option to increase the metadata automatically obtained with a change to the node.cfg. [1]

The metadata is easy to find, and intuitive to use. Files may be found by size, name, or any user defined metadata. See appendix I for examples.

Content can be reindexed at any time using snmp. This versatitilty comes at a price though. All streams must pass through the indexing server so that thier metadata can be indexed. If the indexing server is swamped, or runs out of memory, it can slow the storage cluster to a crawl. Therefor as the cluster increases in size, so to must the number of indexing servers increase.

## 5.6 Caringo Result and Observations

Caringo's CAStor is stateless. This prevents the installation from taking up a disk, but unfortunately this comes at the expense of not being able to customize the storage nodes in any way not specifically supported by the vendor.

A Caringo POSIX interface was provided past the testing deadline. Caringo promises that the Content Filer Server provides a seamless transition into the cluster from a POSIX file system mounted over NFS, and that the CFS can handle multiple writes without causing a bottleneck, but its functionality remains unverified [2].

# 6 Conclusions

We have found that the technology of erasure storage is a viable solution to LANLs archive system due to its scalability, parallelism, and robustness. Individually neither software meets the feature requirements. Caringo has (what appears to be) a more mature solution, but we were unable to test the

POSIX interface. The POSIX interface was one of the requirements for the archive system. Scality has a higher potential to meet the needs of the lab as they are willing to work with their customers to engineer a more fitting solution. The two biggest downfalls of Scality were that they have only beta support for single files over 500MB , and Scality currently requires a stateful (to disk) install. Continued investment into REST interfaced erasure storage has great potential to replace tape drive backups.

## 6.1 Futurework

There is a need to continue testing if erasure storage can meet the speed requirements of high performance computing. There will need to be some calibration done, as the erasure storage systems are currently geared towards cloud (or internet) applications. There are already plans and designated clusters to continue testing and verify the viability of erasure storage.

# 7    Acknowledgements

# 8 References

[1] Caringo. http://www.caringo.com. `http://www.Caringo.com`.

[2] Caringo. Castor scalable, efficient and elastic object storage software. `http://www.Caringo.com/downloads/datasheets/Caringo-CAStor-Object-Storage.pdf`, June 2012.

[3] Gary Grider. ExaScale FSIO and Archive Can we get there? Can we afford to? Technical Report LA-UR-10-04611, Los Alamos National Laboratory, September 2011. Slides for presentation.

[4] J. S. Plank and C. Huang. Tutorial: Erasure coding for storage applications, part 1. Slides presented at FAST-2013: 11th Usenix Conference on File and Storage Technologies `http://web.eecs.utk.edu/~plank/plank/papers/FAST-2013-Tutorial.html`, February 2013.

[5] J. S. Plank and C. Huang. Tutorial: Erasure coding for storage applications, part 2. Slides presented at FAST-2013: 11th Usenix Conference on File and Storage Technologies `http://web.eecs.utk.edu/~plank/plank/papers/FAST-2013-Tutorial.html`, February 2013.

[6] Scality. http://www.docs.scality.com.

[7] Scality. http://www.scality.com. `http://www.Scality.com`.

[8] Marc Staimer. Cloud storage's "organic" or living evolution. Slides presented at 2011 SNIA Cloudburst Summit `http://www.snia.org/sites/default/files2/cloudburst2011/presentations/MarcStaimer_Cloud_Storage_Organic_rev02.pdf`, September 2011.

# Appendices

Appendix A: Network Setup

The Admin subnet was not necessary and can be ignored.

Listing 3: dhcpd.conf

```
1  ddns−update−style interim;
2  option routers 192.168.0.254;
3  allow client−updates;
4  shared−network name{
5      option domain−name "orange.com";
6      option domain−name−servers newHead.orange.com;
7      option subnet−mask 255.255.255.0;
8  #Admin Subnet
9  subnet 192.168.0.0 netmask 255.255.255.0{
10     option subnet−mask 255.255.255.0;
11     option domain−name−servers 192.168.0.254;
12     option routers 192.168.0.254;
13 }
14 #Scality Subnet
15 subnet 192.168.1.0 netmask 255.255.255.0 {
16     option subnet−mask 255.255.255.0;
17     option domain−name−servers 192.168.1.254;
18     option routers 192.168.1.254;
19     authoritative;
20 }
21 #Caringo Subnet
22 subnet 192.168.2.0 netmask 255.255.255.0 {
23     allow booting;
24     allow bootp;
25     filename "pxelinux.0";
26     range 192.168.2.100 192.168.2.200;
27     option subnet−mask 255.255.255.0;
28     option domain−name−servers 192.168.2.254;
29     option routers 192.168.2.254;
30     authoritative;
31 }
32         host caringoadm {
33                 option host−name caringoadm;
34                 option domain−name "orange.com";
35                 hardware ethernet 00:30:48:c5:96:83;
36                 fixed−address 192.168.2.1;
37         }
38 #Scality Nodes
39         host scalityadm {
40                 option host−name scailtyadm;
41                 option domain−name "orange.com";
42                 hardware ethernet 00:30:48:c5:96:73;
43                 fixed−address 192.168.1.1;
44         }
45         host n15 {
46                 option host−name n15;
47                 option domain−name "orange.com";
48                 hardware ethernet 00:26:55:18:92:50;
49                 fixed−address 192.168.1.115;
50         }
51         host n00 {
52                 option host−name n00;
```

```
53          option domain−name "orange.com";
54          hardware ethernet 00:25:b3:ae:cb:cc;
55          fixed−address 192.168.1.100;
56      }
57   host n01 {
58          option host−name n01;
59          option domain−name "orange.com";
60          hardware ethernet 00:26:55:18:82:5e;
61          fixed−address 192.168.1.101;
62      }
63   host n03 {
64          option host−name n03;
65          option domain−name "orange.com";
66          hardware ethernet 00:26:55:18:a5:02;
67          fixed−address 192.168.1.103;
68      }
69   host n10 {
70          option host−name n10;
71          option domain−name "orange.com";
72          hardware ethernet 00:25:b3:ae:86:48;
73          fixed−address 192.168.1.110;
74      }
75   host n07 {
76          option host−name n07;
77          option domain−name "orange.com";
78          hardware ethernet 00:26:55:18:a4:74;
79          fixed−address 192.168.1.107;
80      }
81 }
```

Listing 4: orange.com.fwd

```
1  $ORIGIN orange.com.
2  $TTL 1D
3  @ IN SOA newhead.orange.com. root.orange.com. ( ; SOA is source, so put your dns server name here, and any
        administrators
4                                      0 ; serial
5                                      1D ; refresh
6                                      1H ; retry
7                                      1W ; expire
8                                      3H ) ; minimum
9        IN NS newhead.orange.com. ; NS is loopback device
10 ;Head/Admin subnet
11 newHead IN A 192.168.0.254 ; IN A is a new dns entry for the table
12 scalityadm IN A 192.168.1.1
13 caringoadm IN A 192.168.2.1
14
15 ;Scality subnet
16 n00 IN A 192.168.1.100
17 n01 IN A 192.168.1.101
18 n03 IN A 192.168.1.103
19 n07 IN A 192.168.1.107
20 n10 IN A 192.168.1.110
21 n15 IN A 192.168.1.115
22 scality IN A 192.168.1.100
23 scality IN A 192.168.1.101
24 scality IN A 192.168.1.103
25 scality IN A 192.168.1.107
26 scality IN A 192.168.1.110
27 scality IN A 192.168.1.115
28
29 ;Caringo subnet
30 caringo IN A 192.168.2.102
31 caringo IN A 192.168.2.103
```

```
32 caringo IN A 192.168.2.104
33 caringo IN A 192.168.2.105
34 caringo IN A 192.168.2.106
35 caringo IN A 192.168.2.107
36 caringo IN A 192.168.2.108
37 caringo IN A 192.168.2.109
38 caringo IN A 192.168.2.110
39 caringo IN A 192.168.2.112
40
41 ;IPMI Interfaces
42 ipmi01 IN A 192.168.1.11
43 ipmi02 IN A 192.168.1.12
44 ipmi03 IN A 192.168.1.13
45 ipmi04 IN A 192.168.1.14
46 ipmi05 IN A 192.168.1.15
47 ipmi06 IN A 192.168.1.16
48 ; = comment
49 ; The origin selection will be attached to any address that does not end in a period, so remember to add final .
```

Appendix B: Scality Overview

Each green oval is a virtual node.

Figure 4: Scality Ring

Each virtual node is listed below with its keys.

Figure 5: Scality Keys

| Name | Predecessor | Key | Objects | Used (TB) | Total (TB) | CPU | Tasks | State | Action |
|---|---|---|---|---|---|---|---|---|---|
| n03-n6 | F8E38E | 000000 | 267810 | 1.19 | 2.95 | 3% | 0 | RUN | Leave |
| n00-n6 | 000000 | 071C71 | 261968 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n10-n6 | 071C71 | 0E38E3 | 267610 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n03-n5 | 0E38E3 | 155555 | 267137 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n00-n5 | 155555 | 1C71C7 | 262710 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n01-n6 | 1C71C7 | 238E38 | 262148 | 1.19 | 2.95 | 3% | 0 | RUN | Leave |
| n15-n6 | 238E38 | 2AAAAA | 267823 | 1.19 | 2.95 | 3% | 0 | RUN | Leave |
| n10-n5 | 2AAAAA | 31C71C | 261950 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n15-n5 | 31C71C | 38E38E | 262469 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n00-n4 | 38E38E | 400000 | 267163 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n15-n4 | 400000 | 471C71 | 262646 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n07-n6 | 471C71 | 4E38E3 | 262208 | 1.19 | 2.95 | 3% | 0 | RUN | Leave |
| n00-n3 | 4E38E3 | 555555 | 262720 | 1.19 | 2.95 | 3% | 0 | RUN | Leave |
| n03-n4 | 555555 | 5C71C7 | 267037 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n07-n5 | 5C71C7 | 638E38 | 262536 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n01-n5 | 638E38 | 6AAAAA | 262069 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n03-n3 | 6AAAAA | 71C71C | 267819 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n10-n4 | 71C71C | 78E38E | 262147 | 1.19 | 2.95 | 3% | 0 | RUN | Leave |
| n10-n3 | 78E38E | 800000 | 262713 | 1.19 | 2.95 | 3% | 0 | RUN | Leave |
| n01-n4 | 800000 | 871C71 | 267043 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n00-n2 | 871C71 | 8E38E3 | 262468 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n07-n4 | 8E38E3 | 955555 | 262089 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n10-n2 | 955555 | 9C71C7 | 267771 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n03-n2 | 9C71C7 | A38E38 | 267271 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n07-n3 | A38E38 | AAAAAA | 262715 | 1.19 | 2.95 | 3% | 0 | RUN | Leave |
| n15-n3 | AAAAAA | B1C71C | 261968 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n01-n3 | B1C71C | B8E38E | 262536 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n10-n1 | B8E38E | C00000 | 262068 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n07-n2 | C00000 | C71C71 | 262711 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n00-n1 | C71C71 | CE38E3 | 267219 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n01-n2 | CE38E3 | D55555 | 262714 | 1.19 | 2.95 | 3% | 0 | RUN | Leave |
| n07-n1 | D55555 | DC71C7 | 261950 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n03-n1 | DC71C7 | E38E38 | 267546 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n15-n2 | E38E38 | EAAAAA | 262089 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n01-n1 | EAAAAA | F1C71C | 262645 | 1.19 | 2.95 | 2% | 0 | RUN | Leave |
| n15-n1 | F1C71C | F8E38E | 262208 | 1.19 | 2.95 | 3% | 0 | RUN | Leave |

Page: [1]

Appendix C: Scality Install

## C.1 Node Configuration

Listing 5: Node Config

```
1 yum install −y openssh−clients kpartx parted epel−release
2 /etc/init.d/iptables stop
3 /root/scality/scality−install/scripts/scality−auto−build−disks.sh −X −y −A −i sda
4 cp /root/scality/scality−install/config/centos6/etc/yum.repos.d/scality.repo /etc/yum.repos.d/scality.repo
5 echo "'hostname −i' 'hostname −f'" >> /etc/hosts
6 yum −y install scality−biziod.x86_64 scality−node.x86_64 scality−ringsh.x86_64 scality−sagentd.x86_64 scality−sd−
     httpd.x86_64 scality−sfused.x86_64 scality−sfused−common.x86_64 scality−sproxyd.x86_64 scality−sproxyd−
     httpd.x86_64 scality−srebuildd.x86_64 scality−srebuildd−httpd.x86_64 scality−mod_dewpoint.x86_64 scality−
     ringsh scality−srebuildd
7 scality−node−config −p /scality/disk −d 3 −n 6 −u 192.168.1.1 −m 'hostname −s'−n −i 'hostname −i'
8 scality−sagentd−config −u 192.168.1.1
9 /etc/init.d/scality−node restart
10 /etc/init.d/scality−sagentd restart
11 /usr/local/scality−ringsh/config/generate−ringsh.py https://192.168.1.1:3443 data
```

## C.2 Supervisor Configuration

Listing 6: Supervisor Config

```
1 #the 100,101,103 etc corresponded to the last 3 digits of our ip addresses for each node
2 for i in 100 101 103 104 107 113; do ringsh.py −r data supervisor serverAdd n$i−sagentd 192.168.1.$i 7084 ; done
3 echo −n "bstrap ="; ringsh.py −r data supervisor serverList | awk '{print $7}' | sed s/\:/\ /g | awk '{print $1 "
     :4244"}' | tr '\n' ','
```

## C.3 Node Configuration

Listing 7: sproxyd.conf

```
1
2 #
3 # sproxyd config file
4 #
5
6 #number of worker threads per task pool
7 #n_workers = 100
8
9 #bind address
10 #bind = 0.0.0.0
11
12 #bind port
13 port = 10000
14
15 #connections backlog
16 #backlog = 10000
17
18 #number of incoming requests processed concurrently
19 #n_responders = 500
20
21 #maximum number of open file descriptors
22 max_proc_fd = 40960
23
24 #syslog_facility = "daemon"
```

```
25  #chroot_path =
26  #uid =
27  #gid =
28
29  #max simultaneous connections open to ring nodes
30  conn_max = 10000
31
32  #max number of connection re−use before kicking it
33  conn_max_reuse = 100000
34
35  #enable or disable consistency check on reads, when requested
36  consistent_reads = 1
37
38  #enable or disable consistency check on writes, when requested
39  consistent_writes = 1
40
41  #list of class translations to enable. It's a global parameter that
42  #can't be specified per interface. The list consists of 0 or more
43  #tokens, each token if of the form "a=b" or "a:b" (without the quotes
44  #and without whitespace inside the token), where a is the class stored
45  #in the key, and b is the interpreted class for those keys.
46  #class_translations = "1:2,3:4"
47
48
49  #
50  # base_path=/chord configuration
51  #
52  [chord]
53
54  #
55  # STRUCTURAL PARAMETERS (do not change them when set for a storage Ring)
56  #
57
58  #ring driver
59  ring_driver = "chord"
60
61  #enable path −> key conversion to use arbitrary paths as keys. a SHA1
62  #hash is done on the path after merging any consecutive string of '/' into
63  #a single '/'.
64  by_path_enabled = 0
65
66  #service−id of keys generated by path −> key conversion
67  #by_path_service_id = 0xC0
68
69  #class of service of keys generated by path −> key conversion
70  #by_path_cos = 2
71
72
73  #
74  # RING ACCESS PARAMETERS
75  #
76
77  #bootstraplist (port 4244 by default)
78  #eg. bstraplist = 1.2.3.4:4244,1.2.3.4:4245
79  bstraplist =
        192.168.1.100:4244,192.168.1.101:4244,192.168.1.103:4244,192.168.1.107:4244,192.168.1.110:4244,192.168.1.115:4244,192.168.1.
80
81
82  #
83  # CHORD PERFORMANCE TUNING
84  #
85
```

```
86  #enable asynchronous writes (put) triggered in request headers (see
87  #documentation for how to specify asynchronous behaviour in requests)
88  deferred_writes_enabled_by_request = 0
89
90  #enable asynchronous writes (put) for classes of service defined in
91  #"deferred_writes_policy" parameter
92  deferred_writes_enabled_by_policy = 0
93
94  #set the asynchronous writes policy per class of service. It's
95  #comma−separated or semicolon−separated list of class of service
96  #specs. Each spec is the class of service number, followed by a ':',
97  #followed by the minimum number of successful writes required before
98  #returning a success to the caller.
99  #deferred_writes_policy = "1:1,2:1,3:2,4:2,5:3"
100
101  #enable asynchronous deletes triggered in request headers (see
102  #documentation for how to specify asynchronous behaviour in requests)
103  deferred_deletes_enabled_by_request = 0
104
105  #enable asynchronous deletes for classes of service defined in
106  #"deferred_deletes_policy" parameter
107  deferred_deletes_enabled_by_policy = 0
108
109  #set the asynchronous deletes policy per class of service. It's
110  #comma−separated or semicolon−separated list of class of service
111  #specs. Each spec is the class of service number, followed by a ':',
112  #followed by the minimum number of successful deletes required before
113  #returning a success to the caller.
114  #deferred_deletes_policy = "1:1,2:1,3:2,4:2,5:3"
115
116
117
118  #
119  # base_path=/arc configuration
120  #
121  [arc]
122
123  #
124  # STRUCTURAL PARAMETERS (do not change them when set for a storage Ring)
125  #
126
127  #ring driver
128  ring_driver = "arc"
129
130  #enable path −> key conversion to use arbitrary paths as keys. a SHA1
131  #hash is done on the path after merging any consecutive string of '/' into
132  #a single '/'.
133  #by_path_enabled = 0
134  by_path_enabled = 0
135
136  #service−id of keys generated by path −> key conversion
137  #by_path_service_id = 0xC0
138
139  #class of service of keys generated by path −> key conversion
140  #by_path_cos = 0
141
142
143  #how many equally−spaced areas the RING is configured to support
144  arc_schema = 6
145
146
147  #
148  # RING ACCESS PARAMETERS
```

```
149  #
150
151  #bootstraplist (port 4244 by default)
152  #eg. bstraplist = 1.2.3.4:4244,1.2.3.4:4245
153  bstraplist =
           192.168.1.100:4244,192.168.1.101:4244,192.168.1.103:4244,192.168.1.107:4244,192.168.1.110:4244,192.168.1.115:4244,192.168.1.
154
155
156  #
157  # ARC STORAGE OVERHEAD CONTROL
158  #
159
160  #object class to use for replication
161  object_class = 2
162
163  #minimum size in bytes where object is not stored replicated
164  #(−1: always replicated)
165  replication_size_threshold = 60000
166
167  #maximum number of arc data chunks
168  n_data_parts = 3
169
170  #number of arc coding chunks
171  n_coding_parts = 3
172
173  #number of redundant chunks written successfully before a PUT
174  #is considered successful (−1: all writes must succeed)
175  min_redundant_parts_put_ok = 2
176
177  #minimum size in bytes of a arc data chunk
178  min_data_part_length = 20000
179
180
181  #
182  # ARC PERFORMANCE TUNING
183  #
184
185  #number of parallel tasks in the main task pool
186  main_n_workers = 10
187
188  #number of parallel tasks in the sub task pool
189  sub_n_workers = 10
190
191  #number of parallel tasks in the cache task pool
192  cache_n_workers = 10
193
194  #size in bytes of data stripes when cutting data on put
195  #stripe_size = 262144
196  stripe_size = 1048576
197
198  #size in bytes of data buffers read at once when reconstructing
199  get_reconstruct_buffer_size = 1048576
200
201  #maximum number of write buffers queued from each client before blocking
202  max_stripes_in_write_queue = 2
203
204  #enable or disable caching of data and metadata
205  #(speeds up gets that need reconstruction)
206  chordcache_enabled = 1
207
208
209
```

```
210 #sproxyd HOWTO using curl:
211
212 #do a PUT with base64−encoded user metadata ("myusermd")
213 #curl −0 −XPUT −H "x−scal−usermd: bXl1c2VybWQ=" http://localhost:81/proxy/chord/88
        FF0A8375F3112C8E340A38E38FE93438412120 −−data−binary @/etc/hosts
214
215 #do a UPDATEMD (put partial): ("newusermd")
216 #curl −0 −XPUT −H "x−scal−cmd: update−usermd" −H "x−scal−usermd: bmV3dXNlcm1k" http://localhost
        :81/proxy/chord/88FF0A8375F3112C8E340A38E38FE93438412120
217
218 #do a GET
219 #curl http://localhost:81/proxy/chord/88FF0A8375F3112C8E340A38E38FE93438412120
220
221 #do a GET with a range (retrieve 500 bytes from the byte 1000).
222 #ranges with at least one missing bound are not supported.
223 #curl −r 1000−1499 http://localhost:81/proxy/chord/88FF0A8375F3112C8E340A38E38FE93438412120
224
225 #do a STAT (dumps all HTTP headers received, user metadata included)
226 #curl −I http://localhost:81/proxy/chord/88FF0A8375F3112C8E340A38E38FE93438412120
227
228 #do a DELETE
229 #curl −XDELETE http://localhost:81/proxy/chord/88FF0A8375F3112C8E340A38E38FE93438412120
230
231
232 #get statistics information
233 #curl http://localhost:81/proxy/chord/.stats
234
235 #reset statistics information
236 #curl −XDELETE http://localhost:81/proxy/chord/.stats
237
238 #get active configuration parameters, each on a line with format key=value
239 #curl http://localhost:81/proxy/chord/.conf
```

## C.4  Supervisor Configuration

Listing 8: sfused.conf

```
1  #
2  ## Scality sfused configuration file
3  ##
4  ##
5
6  ## Documentation is available here:
7  ## http://docs.scality.com/doku.php?id=sfused:bizfs_overview#sfused
8
9  # Volume number (0..2^32)
10 # Scoping: when you change the 'dev' value, you don't see the objects put with another device id
11 dev = 987654323
12
13 # Dynamic mountpoint
14 mountpoint = /ring/fuse
15
16 # The logger_id string is added to the log lines to mark the logs concerning this conf.
17 logger_id = sfused
18
19 # check /usr/include/sys/syslog.h to see all the possibilities
20 # e.g.: auth, authpriv, news, syslog, mail, etc
21 syslog_facility = mail
22
23 conn_max = 10000
24 conn_tcp_nodelay = 1
25 conn_max_reuse = 100000
26
```

```
27  # check that a user belongs to a posix group (from /etc/group or whatever)
28  # note that it might degrade the performance because checks are performed very often
29  group_check = 0
30
31  ##
32  ## Ring FS
33  ##
34
35  file_cos = 5
36  cat_cos = 5
37  cat_page_cos = 5
38  dir_cos = 5
39  dir_page_cos = 5
40  rootfs_cos = 5
41  n_open_files_buckets = 65521
42  honor_forget = 1
43  max_proc_fd = 40960
44
45
46  ##
47  ## RootFS
48  ##
49
50  rootfs = 1
51  rootfs_type = md5
52  rootfs_cache = 3
53  allow_rootfs_listing = 1
54
55
56  ##
57  ## Cache
58  ##
59
60  # Set this property to allow a nonprivileged user to create/remove/rename directories under rootfs
61  allowed_rootfs_uid = 501
62
63  # Disable version checking in cache, semi−stateless mode. (recommended in failover setting)
64  cache_enable_checks = 0
65
66  # Force a check every x seconds for object read access, default is 0 (check all the time)
67  cache_check_time = 0
68
69  # Files preview (for application reading mail headers)
70  # To allow file previews set cache_preview_enable to "1", recommended value is disable, '0' (default value)
71  cache_preview_enable = 1
72  cache_preview_bytes = 16384 # in bytes
73
74  cache_preview_tail_trick = 1
75  cache_serialization_dir = /var/cache/sfused_cache
76  cache_serialization_period = 600
77
78  ##
79  ## Pools of worker threads
80  ##
81
82  # Chord workers
83  n_workers = 480
84
85  # workers_arc_main =
86  # workers_arc_sub =
87  # workers_arc_cache =
88  # workers_db_index =
89  # workers_db_commit =
```

```
 90  # workers_db_delete =
 91  # workers_prefetch =
 92
 93  ##
 94  ## Undelete
 95  ##
 96
 97  # Enable undelete journal to keep track of deleted keys
 98  undelete = 1
 99
100  # Purge threshold, in seconds, for undelete journal
101  undelete_purge_threshold = 604800
102
103  # Max age, in seconds, of deleted files in journal. Default is 604800 seconds (6 days)
104  # Value must be lower than the Ring purge expiration time
105  undelete_duration = 300
106
107  # Allow undelete for objects matching this regexp (u\.* = dovecot mails)
108  undelete_pattern = "^u[.].*"
109
110  # Do we match only the object's name or its full path?
111  undelete_pattern_full_path = 0
112
113  ##
114  ## geosync
115  ##
116
117  #geosync = 1
118  #geosync_interval = 30
119  #geosync_prog = /usr/local/bin/ssync
120  #geosync_args = 'ssync $FILE $MOUNTPOINT / --dev $DEV --file-cos $FILE_COS --dir-cos $DIR_COS'
121
122
123  [fuse]
124  fuse_max_tasks = 72
125  direct_io = 0
126  big_writes = 1
127  read_ahead_kb = 128
128
129  [ring_driver:0]
130  type = chord
131  #bstraplist = node01-fuse.scality.com,node02-fuse.scality.com,node03-fuse.scality.com
132  bstraplist =
         192.168.1.100:4244,192.168.1.101:4244,192.168.1.103:4244,192.168.1.107:4244,192.168.1.110:4244,192.168.1.115:4244,192.168.1.
133
134  client_routing = 0
135
136  [ring_driver:1]
137  type = srest
138
139  # The base_path needs to match the path to the correct driver on sproxyd.
140  base_path=/proxy/arc/
141  #
142  # ## <CHANGE_ME>
143  #
144  # # NOTE: The data ring should be compromised of the first 6 node daemon processes on each storage node.
145  # # In this example, node daemon 1 on the first 6 servers is chosen (the port, if not listed, defaults to 4244).
146  # # The port used is port 81 for the httpd FastCGI module supporting sproxyd.
147  bstraplist =
         192.168.1.100:81,192.168.1.101:81,192.168.1.103:81,192.168.1.107:81,192.168.1.110:81,192.168.1.115:81
148  #
149  #
```

```
#Files cache

[cache:0]
ring_driver = 0
size = 2000000000
serialization = 1


#Directories cache

[cache:1]
ring_driver = 0
size = 2000000000
serialization = 1

[cache:2]
ring_driver = 0
size = 2000000000
serialization = 1

[cache:3]
ring_driver = 0
size = 100000000
serialization = 1

#File Inode Mode

[ino_mode:0]
type = mem
max_file_size = 536870912
cache = 0

#Directories Inode Mode

[ino_mode:3]
type = mem
cache = 1

#Asynchronous Inode Mode

[ino_mode:1]
type = async
max_file_size = 536870912
cache = 2
pattern = 'dovecot.index|.temp.centos6−fuse−pds.scality.com|.deleted'
queue_path = /var/cache/sfused
pattern_full_path = 0
period = 30
ttl = 60
n_workers = 12
size = 4000000000

##Sparse Inode Mode
#
[ino_mode:2]
type = sparse
cache_stripes = 0
cache_md = 0
pattern = .*
pattern_full_path = 0
sticky = 1
```

```
213 #
214 fsid = 1
215 main_cos = 5
216 page_cos = 5
217 stripe_cos = 5
218 stripe_size = 1048576
219 #
220 file_dirty_limit = 128000000 # 128MB
221 global_dirty_limit = 1000000000 # 1GB
222 dirty_timeout = 5
223 fsync_on_close = 1
224 workers_io = 64
225 workers_commit = 64
```

Appendix D: Scality cURL Commands

Listing 9: Scality cURL

```
1  #Verify "scality−sproxyd" Service is Running
2  #NOTE: This service needs to be running on the ring nodes, if not running then start it.
3  #sudo su
4  #service scality−sproxyd status
5
6  #Generate a Random Key
7  ringsh.py
8
9  #When you are in the ringsh.py program enter the following command:
10 key random
11 #Use the key that it prints to the screen in your get/put/delete curl commands.
12
13 #PUT a File
14 curl −XPUT −H "Expect:" −H "x−scal−usermd: bXl1c2VybWQ=" http://localhost:81/proxy/chord/70
       A249E993315C2CC3A7F10F76BCC37099447B00 −−data−binary @/root/test
15
16 #GET a File
17 curl http://localhost:81/proxy/chord/70A249E993315C2CC3A7F10F76BCC37099447B00
18
19 #DELETE a File
20 curl −XDELETE http://localhost:81/proxy/chord/70A249E993315C2CC3A7F10F76BCC37099447B00
21
22
23 ##ARC
24 #Generate a random ARC key, the last bit of the key is not random, but specifies settings
25 scalkeyarcgen −t arc −k 3 −m 3 −s 2 UPLOAD_FILE_PATH/FILE_NAME
26 Output: D41D8CD98F00B204E9800900000000510C302070
27
28 #PUT a File
29 curl −0 −XPUT −H "x−scal−usermd: bXl1c2VybWQ=" http://localhost:81/proxyd/arc/
       D41D8CD98F00B204E9800900000000510C302070 −−data−binary @UPLOAD_FILE_PATH/FILE_NAME
30
31 #Get a File
32 curl http://localhost:81/proxyd/arc/D41D8CD98F00B204E9800900000000510C302070 >
       DOWNLOAD_FILE_PATH/FILE_NAME
33
34 #Delete a File
35 curl −−XDELETE http://localhost:81/proxyd/arc/D41D8CD98F00B204E9800900000000510C302070
36
37 ##CDMI (We didn't have these commands functioning)
```

```
38 curl −0 −v −XPUT −H "Content−Type: application/cdmi−object" −H "X−CDMI−Specification−Version: 1.0.1"
      −d '{ "metadata": { "Data":"bXl1c2VybWQ" }, "value": "some text" }' http://jumphost:80/fs/chris80
```

## Appendix E: Scality MetaData

### Listing 10: ListKeys.sh

```bash
1  #!/bin/bash
2  DSO=$1 #data
3  listKeysMethod=$2 #index or browse
4  db=$3
5  table=$4
6  export TERM=vt100
7  nodes=$(ringsh.py supervisor dsoStatus $DSO | grep "Node:" | cut −d" " −f 2 )
8  now=$(date +%s)
9  for node in $nodes ; do
10     echo "Listing keys on $node"
11     echo −e "load conf $DSO \n\n use $node \n\n node listKeys loadmetadata=$listKeysMethod \n\n" | ringsh.py
          > work/keys.$node.$now
12  # I just redirect them from the above echo to a file, not an sql database (as the below lines would do)
13  # echo "generating sql file for loading data from $node"
14  # sed "s/\,/\'\,\'/g" work/keys.$node.$now | sed "s/$/\'\)\;/g" | sed "s/^/insert into $table values\(\'$node
        \'\,\'$listKeysMethod\'\,\'/" > work/keys.$node.$now.sql
15  # echo ".quit" >> work/keys.$node.$now.sql
16  # echo "loading data from $node in $db $table"
17  # cat work/keys.$node.$now.sql | sqlite3 $db
18  done
19  # xargs −t −n 1 −I{} echo −e \"load conf ring\n\n use {}\n\nnode dumpStats
```

## Appendix F: Scality SFUSE

### Listing 11: SFUSE Install

```bash
1  #Scality Documentation:
2  #http://docs.scality.com/display/DOCS/Install+the+SOFS+Connector+on+CentOS+or+RedHat
3
4  #Install Scality Package
5  sudo su
6  yum install scality−mod_dewpoint−4.1.3.r34599−1.el6.x86_64
7  #This may require a scality repository if I remember correctly:
8  #[scality−base]
9  #name=CentOS6 − Scality Base
10 #baseurl=http://scalitycs:CSrepoPass@packages.scality.com/stable_isildur/centos/6/x86_64/
11 #gpgcheck=0
12
13 #Copy Scality Configuration Files
14 cp /scality/scalityinstall/configs/scality/etc/ /etc/
15 scp nXX:/etc/sfused.conf /etc/sfused.conf
16 #Note: nXX is the DNS name or IP address of a storage node in the Scality ring.
17
18 #Create "/ring" directory and Set Permissions
19 mkdir /ring
20 chmod g+w /ring
21 chown root:fuse /ring
22
23 #Mount the Control Filesystem for Fuse
```

```
24 mount −t fusectl none /sys/fs/fuse/connections
25
26 #Disable Transparent Hugepage Option
27 echo never > /sys/kernel/mm/redhat_transparent_hugepage/enabled
28
29 #Install Connector
30 yum install scality−sfused
31
32 #Create a Catalogue
33 sfused −X −c /etc/sfused.conf
34
35 #Start Connector
36 /etc/init.d/scality−sfused start
37
38 #Test Mounted Volume
39 mount −l | grep fuse
```

Appendix G: Caringo Install

Listing 12: Caringo Config Files

```
1  #################### Sample pxelinux.cfg ################
2  timeout 100
3  default menu.c32
4  ONTIMEOUT 1
5
6  menu title ############## CASTor boot menu ###########
7  label 1
8       kernel profiles/castor/kernel
9       append initrd=profiles/castor/fsimage ramdisk_size=128000 root=/dev/ram0 castor_cfg=http://192.168.2.1/
            castor/cfg−list.txt
10
11 ############## Sample cfg−list.txt #############
12 http://192.168.2.1/castor/node.cfg
13
14 ############## Sample license.txt ##############
15 licenseFormat = 1.1
16
17 cn = Los Alamos National Laboratory
18 street =
19 l = Los Alamos
20 st = NM
21 postalCode =
22 co = US
23
24 clusterDescription = 1st Test Bed
25
26 # License Components
27 expirationDate = 2013−09−01
28 featureClusterMaxTB = 50
29 featureContentIndexing = yes
30 featureErasureCoding = yes
31 featureMinimumMinReps = 1
32
33 −−−−−BEGIN PGP SIGNATURE−−−−−
34 Version: GnuPG v1.4.10 (GNU/Linux)
35
36 iEYEARECAAYFAlHTCbEACgkQRYikRJU1RfOrvACeLkPs9Ro5A7HKxzYc4qfyz01p
37 +xMAn1M8h0YbKl3QUSCRDz4fZwuWf6ND
38 =Vdv2
```

```
39  -----END PGP SIGNATURE-----

40

41  ############### Sample Node configuration file ################

42

43  #to set a value the syntax is [key] = [value]
44  #[key] = type of setting to change
45  #[value] = the value to change the setting to
46  #a one space seperation is expected between the [key] '=' and [value]
47  #All headers can be chained together with dot notation.
48  #E.G. [cluster] \n\t name is equivalent to cluster.name

49

50

51  ########    Cluster Name
52  [cluster]
53          name = caringo.orange.com

54

55

56  ########  Node Disk Formating
57  #Sets the nodes to use all disk volume availible
58  #the :k flag sets the volumes to be kept if they expire.
59  #Volumes are set to expire if they are not detected in the last 14 days
60  #the cluster was active
61  [disk]
62          volumes = all:K

63

64

65  ########   NTP TIME SOURCE
66  #sets the timesource to the local ntp server
67  [network]
68          timeSource = 192.168.2.254
69  ########    LOG SERVER SETUP

70

71  #Sets log host ip, the minimum error level to log, and posts object id's
72  #to the log file if that object has an error, instead of hiding object id's
73  [log]
74          host = 102.168.2.1
75          level = 40
76          obscureUUIDs = False

77

78  ########   License SETUP
79  #Castor requires that the product licence be serviced by a web server
80  [license]
81          url=http://192.168.2.1/castor/license.txt

82

83

84  ########   Security SETUP
85  #You may also set many security options. For our test bed we left the snmp public
86  #format for users is {'USER' : 'PASSWORD' }
87  #Instead of clear text use MD5 of following string '<username>:CAStor administrator:<password>'
88  [security]
89          administrators = {'admin' : 'admin' , 'snmp' : 'admin'}
90          operators = {'snmp':'public'}

91

92  ########  ERASURE CODE SETUP
93  #To set the cluster to erasure code add the following lines to node.cfg
94  #sets encoding to 3 + 3 and minimume file size to erasure code to 60KB
95  [ec]
96          encoding = 3:3
97          minStreamSize = 61440

98

99

100

101
```

```
102  ######## INDEXER SETUP
103  #The only supported indexer service is elastic search, using the
104  #comercial plugin provided by caringo. This set up sets elastic search
105  #as the indexing agent, and sets the metadata kept on each stream in the
106  #cluster to full
107  [indexer]
108          name = elasticsearch
109          fullMetadata = 1
110
111  [indexer_elasticsearch]
112          host = HOST_IP
113          port = 9200 #Default port number, can be changed
114          connectionRetryInterval = 10
115          insertBatchSize = 100
116          insertBatchTimeout = 60
117          className = caringo.castor.indexer.plugin.elasticsearchplugin.ElasticSearchPlugin
```

Appendix H: Caringo cURL Commands Cluster name = caringo.orange.com Domain name = caringo.orange.com (Doesn't have to be the same as cluster name )

Listing 13: Caringo Basic Commands

```
1  #BUCKET CREATION EXAMPLE
2  curl −i −−post301 −−data−binary '' −−location−trusted 'http://caringo.orange.com/mybucket?domain=caringo.
      orange.com'
3
4  #PUSH example 'push up hello world'
5  curl −i −−post301 −−data−binary '<html><h1>Hello World</h1></html>' −H 'Content−type: text/html'
      −−location−trusted 'http://caringo.orange.com/mybucket/helloworld?domain=caringo.orange.com'
6
7  #GET example 'get hello world'
8  curl −i −−post301 −−location−trusted 'http://caringo.orange.com/mybucket/helloworld
9  #GET example GUID
10 curl −i −−post301 −−location−trusted 'http://caringo.orange.com/[GUID GOES HERE]?domain=caringo.orange.
      com
11
12 #COPY example 'rename a bucket'
13 #change mybucket to bucket
14 curl −i −−post301 −X COPY −−data−binary '' −−location−trusted 'http://caringo.orange.com/mybucket?
      domain=caringo.orange.com&newname=bucket'
15
16 #Delete a file
17 curl −i −−post301 −X DELETE −−location−trusted 'http://caringo.orange.com/mybucket/helloworld'
```

Appendix I: Caringo MetaData

Listing 14: Caringo Metadata

```
1  #Basic metadata
2  Query arg   #Description
3  tmBorn     #time of create or last update
4  size    #size in bytes
5  name    #UUID or name using URL encoding
6  content−type #content type
7  etag     #entity tag
```

```
 8  sizewithreps  #number of bytes using the maximum reps value
 9
10  #Full metadata
11  content−base
12  content−disposition
13  content−encoding
14  conten−language
15  content−length
16  content−location
17  content−md5
18  last−modified
19  lifepoint
20
21
22  #Custom Fields
23  x−∗−meta
24  x−∗−meta−∗
```

The custom fields can be filled by adding to the http header, and then queried using Get requests. Examples

Listing 15: Metadata Get Requests

```
1  #Query a cluster to find all the the buckets in a domain
2  curl −i http://caringo.orange.com?domain=caringo.orange.com&format=json
3
4  #Query a bucket for all of its contents
5  curl −i http://caringo.orange.com/mybucket?domain=caringo.orange.com&format=xml
6
7  #Find the amount of storage being used by a bucket
8  curl −i http://caringo.orange.com/mybucket?domain=caringo.orange.com&format=json&du=yes
```

# Glossary

**Block storage** This is the underlying idea behind most file systems. The space available is blocked out to store files and folders on. 3, 9

**curl** Linux command used to do http request via the command line. 'man curl'. 6, 8

**FUSE** Filesystem in USerspacE. Tool to enabled Unix operating systems to create a file system in user space. A type of "bridge" to the kernel calls. 6

**Object storage** Uses flexible data containers that are not set to a particular block size. The data is stored uninterpreted with metadata. 3

**POSIX** Portable Operating System Interface. IEEE specifications for maintaining compatibility across multiple operating system.. 3, 6, 8–10

**REST** REpresentational State Transfer. Makes use of the HTTP methods GET, POST, PUT and DELETE to execute operations over HTTP. This is in contrast to SOAP. 5, 6, 8, 10

**SFUSE** Scality's version of FUSE to connect a user to a seeming Block storage device that is actually their ring or Object storage storage in the background. 3, 5, 6